

Lecture 3: Modern Methods of Statistical
Learning sf2935 :
Neural Networks: Learning
Timo Koski

TK

2018-09-03



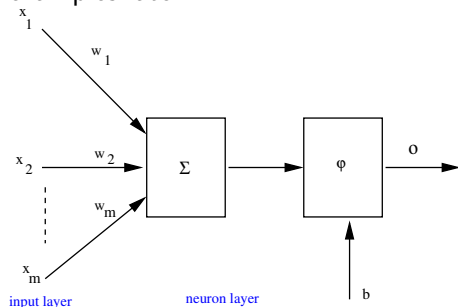
KTH Matematik

- Forward neural networks with multiple layers
- The supervised learning task for neural networks
- Backpropagating
- Neural networks and statistics
- Bias-variance trade-off



Neuron: Input Layer and One Neuron Layer

$\mathbf{x} = (x_1, \dots, x_m) \in \mathbb{R}^m =$ input space, an inner product space of finite dimension $= m$. $\varphi(x) =$ activation function or activator, examples later.



$$\mathbf{w} = (w_1, \dots, w_m)$$

$$o = \varphi \left(\sum_{i=1}^m w_i x_i + b_j \right)$$

\mathbf{w} is called the weight vector and b is a bias.



Sigmoidal activation function

$\varphi(\cdot)$ is a sigmoidal function, if it is nonconstant, bounded, and monotonically-increasing continuous function.

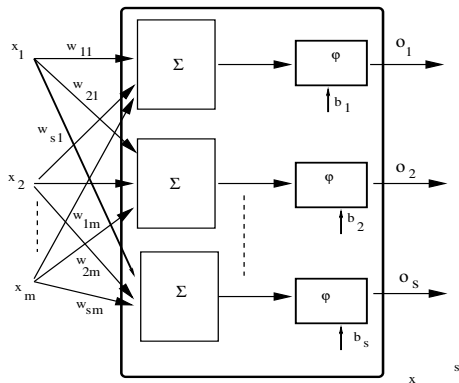
Example

The logistic function:

$$\varphi(z) = \frac{1}{1 + e^{-z}}$$



Neural Network: Input Layer & a Layer of s Parallel Neurons



$$\mathbf{w}_j = (w_{j1}, \dots, w_{jm}) \quad j = 1, \dots, s$$
$$O_j = \varphi \left(\sum_{i=1}^m w_{ji} x_i + b_j \right), \quad j = 1, \dots, s$$



Composition of Functions

$f(x)$ and $g(x)$ are two functions such that the domain of f includes the range of g , the composition $(f \circ g)(x)$ of $f(x)$ and $g(x)$ is defined as

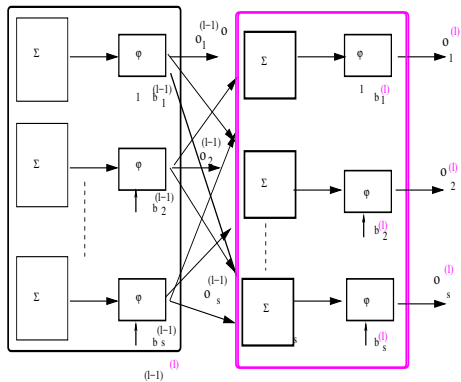
$$(f \circ g)(x) \stackrel{\text{def}}{=} f(g(x)).$$

Chain rule of differentiation:

$$\frac{d}{dx}(f \circ g)(x) = f'(g(x)) \cdot g'(x).$$



Composition of Layers



Output of layer l $o_j^{(l)} = \varphi \left(\sum_{i=1}^s w_{ji}^{(l)} o_i^{(l-1)} + b_j^{(l)} \right), j = 1, \dots, s$

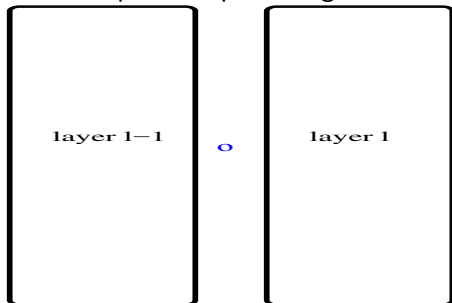
$b_j^{(l)}$ = bias of the j th neuron in the l th layer

$w_{ji}^{(l)}$ = the weight of the i th neuron in the l th layer to the j th neuron in the l th layer.

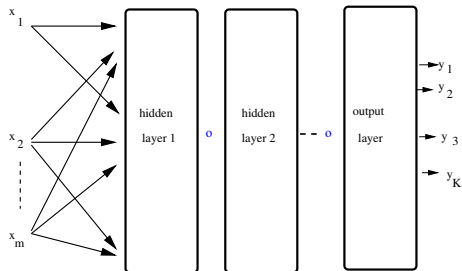
$o_i^{(l-1)}$ = the output of the i th neuron in the $l - 1$ th layer.

Composition of Layers

Let us depict the preceding as



Neural Network with L layers (composed), layer number L is the output layer



Additional layers composed

$y_j = \varphi \left(\sum_{i=1}^{m_{L-1}} w_{ji}^{(L)} o_i^{(L-1)} + b_j^{(L)} \right)$, $j = 1, \dots, K$, the outputs of the neural network. The layers preceding the output layer are called hidden layers, as we do not register their outputs. There is **no feedback**, all computation is feedforward from input layer to output layer.



KTH Matematik

t is the target of the input \mathbf{x} .

Training set: $\mathcal{S} = \left\{ \left(\mathbf{x}^{(r)}, t^{(r)} \right) \right\}_{r=1}^n$ consists of the n data pairs of inputs labelled with the targets $t^{(r)}$.



Examples of target domains

- Supervised classification with K classes: target space consists of vectors t with $K - 1$ zeros and a single one. Hence

$$t = (0, 0, \dots, \underbrace{1}_{\text{k:th position}}, 0, \dots, 0)$$

indicates class k ,

- Estimation of an unknown real valued function f from n values. $K = 1$,

$$t^{(r)} = f(\mathbf{x}^{(r)}), r = 1, 2, \dots, n,$$

the target space = \mathbb{R} .



Quadratic cost function

Since the outputs of the neural network are

$$y_j = \varphi \left(\sum_{i=1}^{m_{L-1}} w_{ji}^{(L)} o_i^{(L-1)} + b_j^{(L)} \right), j = 1, \dots, K$$

we can write them as (composed) functions of the inputs

$$y_j = y_j(\mathbf{x}).$$

We set

$$C(\mathbf{w}, \mathbf{b}) \stackrel{\text{def}}{=} \frac{1}{2n} \sum_{r=1}^n \sum_{i=1}^K (y_i(\mathbf{x}^{(r)}) - t_i^{(r)})^2$$

\mathbf{w} denotes the collection of all the weights in the network, \mathbf{b} all the biases, n is the total number of training inputs,

$\mathbf{y} = (y_1(\mathbf{x}^{(r)}), \dots, y_L(\mathbf{x}^{(r)}))$ is the vector of outputs from the network when $\mathbf{x}^{(r)}$ is input, and the sum is over all training inputs.



Learning algorithms which can automatically tune the weights and biases of a network of artificial neurons so that

$$C(\mathbf{w}, \mathbf{b}) = \frac{1}{2n} \sum_{r=1}^n \sum_{i=1}^K (y_i(\mathbf{x}^{(r)}) - t_i^{(r)})^2$$

is minimized, where $\left\{ \left(\mathbf{x}^{(r)}, t^{(r)} \right) \right\}_{r=1}^n$ is the training set. This is called Minimization of the Mean Squared Error (MSE).



Gradient descent is an iterative algorithm for finding the minimum of a differentiable function. To find a local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient of the function at the current point. Or, in the neighborhood of a point, a function decreases fastest if one goes from the point in the direction of the negative gradient of the function.



Gradient Descent

We write $C = C(\mathbf{w}, \mathbf{b})$. The gradient descent method involves in this case calculating the gradient of C with respect to the weights of the network.

We have a current value of a weight w_{ij}^{old} (dropping the layer index (l) for simplicity of writing). We update the w_{ij}^{old} using gradient descent,

$$w_{ij}^{old} \rightarrow w_{ij}^{old} + \Delta w_{ij}.$$

The change in weight, which is added to the old weight, is equal to the product of the learning rate η and the gradient, multiplied by -1 :

$$\Delta w_{ij} = -\eta \frac{\partial C}{\partial w_{ij}}.$$

The -1 is required in order to update in the direction of a minimum, not a maximum, of the cost function.

Backpropagation means computing the partial derivatives $\frac{\partial C}{\partial w_{ij}^{(l)}}$ and $\frac{\partial C}{\partial b_j}$ in a recursive manner from the output layer to the input



Backpropagation actually allows us to compute the partial derivatives $\frac{\partial C}{\partial w_{ij}}$ and $\frac{\partial C}{\partial b_j}$ for a single training example. We can then recover $\frac{\partial C}{\partial w_{ij}}$ and $\frac{\partial C}{\partial b_j}$ by averaging over training examples. In fact,

- we fix a single training example \mathbf{x} . We remove it from notation as it is a notational nuisance.
-

$$C = \frac{1}{2} \sum_{i=1}^K (t_i - y_i(\mathbf{x}))^2.$$

The factor $\frac{1}{2}$ cancels the exponent, when differentiating.



For simplicity of writing we drop also the layer index (l) in

$$o_j = \varphi \left(\sum_{i=1}^m w_{ji} o_i + b_j \right).$$

We set

$$\text{net}_j \stackrel{\text{def}}{=} \sum_{i=1}^m w_{ji} o_i + b_j,$$

so that

$$o_j = \varphi(\text{net}_j).$$

The input net_j to an activator φ is the weighted sum of outputs o_k of neurons in the previous layer.

If the neuron is in the first layer after the input layer, the o_k of the input layer are simply the inputs x_k to the network. The number of input units to the neuron is m . The variable w_{ji} denotes the weight between neuron i at previous layer and neuron j at the layer above.



Calculation of the partial derivative of the error C with respect to a weight w_{ji} is done using the chain rule twice:

$$\frac{\partial C}{\partial w_{ji}} = \frac{\partial C}{\partial o_j} \frac{\partial o_j}{\partial w_{ji}} = \frac{\partial C}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}} \quad (1)$$



$$\frac{\partial C}{\partial w_{ji}} = \frac{\partial C}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}}$$

The logistic function

$$\varphi(z) = \frac{1}{1 + e^{-z}}$$

has the derivative:

$$\frac{d\varphi(z)}{dz} = \varphi(z)(1 - \varphi(z))$$

Assume that the logistic function is used, then the middle factor in the right hand side of (1) is

$$\begin{aligned} \frac{\partial o_j}{\partial \text{net}_j} &= \frac{\partial}{\partial \text{net}_j} \varphi(\text{net}_j) \\ &= \varphi(\text{net}_j)(1 - \varphi(\text{net}_j)). \end{aligned}$$



$\frac{\partial C}{\partial w_{ji}} = \frac{\partial C}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}}$: Output Layer

The first factor $\frac{\partial C}{\partial o_j}$ in the right hand side of (1), the derivative of the output of $y_j = o_j$ with respect to its input is simple, if the neuron is in the output layer, as then $y_j = o_j$ and

$$\begin{aligned} \frac{\partial C}{\partial o_j} &= \frac{\partial C}{\partial y_j} \\ &= \frac{\partial}{\partial y_j} \frac{1}{2} \sum_{i=1}^K (t_i - y_i)^2 = y_j - t_j = o_j - t_j. \end{aligned}$$



In the last factor of in the right hand side of (1), only one term in net_j depends on w_{ji} . Therefore

$$\frac{\partial \text{net}_j}{\partial w_{ji}} = \frac{\partial}{\partial w_{ji}} \left(\sum_{k=1}^n w_{jk} o_k \right) = \frac{\partial}{\partial w_{ji}} w_{ji} o_i = o_i.$$

Output Layer: Put it together

For the weights in the output layer we have now obtained in the above for (1):

$$\begin{aligned}\frac{\partial C}{\partial w_{ji}} &= \frac{\partial C}{\partial o_j} \frac{\partial o_j}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}} \\ &= (o_j - t_j) \cdot o_j(1 - o_j) \cdot o_i\end{aligned}$$



To calculate here, we first introduce an intermediate quantity, δ_j^l , which we call **the error in the j th neuron (in the l th layer)**. We drop l for simplicity, and define

$$\delta_j \stackrel{\text{def}}{=} \frac{\partial C}{\partial \text{net}_j} \quad (2)$$

We have by the chain rule of calculus

$$\frac{\partial C}{\partial w_{ji}} = \frac{\partial C}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}} = o_j \delta_j. \quad (3)$$

and

$$\delta_j = \frac{\partial C}{\partial \text{net}_j} = \frac{\partial C}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} = \varphi'(\text{net}_j) \frac{\partial C}{\partial o_j}. \quad (4)$$



$\frac{\partial C}{\partial w_{ji}}$: the difficulties at earlier layers

In the calculations to follow the difficulty is that one must make very clear what is fixed. Much (?) of the machine learning literature does not do this with backpropagation and is well confused or difficult to understand.

When computing $\frac{\partial C}{\partial w_{ji}}$ we regard C as function of all the weights. Changes in w_{ji} affect the net _{j} and o_j and all neurons taking input from neuron j including some associated os.



When we take partial derivatives w.r.t. to a net or an output of a neuron, we allow all other signals in the network which depend on the net or the output of a neuron to follow their usual dependence. Thus all weights and all nets (and hence corresponding o outputs) of neurons **in the same and earlier layers are kept fixed.**



We evaluate $\delta_j = \frac{\partial C}{\partial \text{net}_j}$ by noting that net_j affects the neuron outputs through o_j , and **this only acts through connections to other neurons's outputs at layers closer to the network output layer.**



By (4) we have

$$\delta_j = \varphi'(\text{net}_j) \frac{\partial C}{\partial o_j}.$$

Let now o_k be a neuron at the next layer ($l + 1$, this is closer to the output layer) and taking input from neuron net_j at layer l . We indicate this with

$$k : j \rightarrow k$$



Then we take the total derivative w.r.t. to o_j to obtain

$$\frac{\partial C}{\partial o_j} = \sum_{k:j \rightarrow k} \frac{\partial C}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial o_j}$$

and by chain rule

$$= \sum_{k:j \rightarrow k} \frac{\partial C}{\partial o_k} \frac{\partial o_k}{\partial \text{net}_j} w_{kj}$$

and by the middle equality in (4)

$$= \sum_{k:j \rightarrow k} \delta_k w_{kj}$$

$$\frac{\partial C}{\partial o_j} = \sum_{k:j \rightarrow k} \delta_k w_{kj}$$

and

$$\delta_j = \varphi'(\text{net}_j) \sum_{k:j \rightarrow k} \delta_k w_{kj}$$

Therefore, the derivative with respect to o_j can be calculated if all the derivatives with respect to the outputs o_k of the next layer - the one closer to the output neuron - are known.



Put it all together:

$$\frac{\partial C}{\partial w_{ij}} = o_i \delta_j$$

where

$$\delta_j = \frac{\partial C}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} = \begin{cases} (o_j - t_j) o_j (1 - o_j) & \text{if } j \text{ is an output neuron,} \\ (\sum_{k:j \rightarrow k} \delta_k w_{kj}) o_j (1 - o_j) & \text{if } j \text{ is an inner neuron.} \end{cases}$$



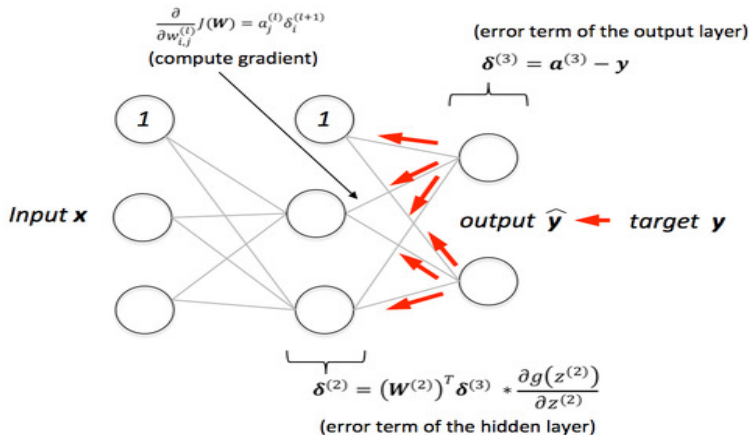
Backpropagation:

The error vectors δ_j are computed backward, starting from the final layer. The backward movement is a consequence of the fact that the cost is a function of outputs from the network.

- a **forward pass** to calculate the outputs from the inputs
- a **backward pass** to calculate δ_j



Backpropagation:



Who discovered backpropagation?



KTH Matematik



Neural Network: Universal Approximator (Cybenko & Hornik)

Let $\varphi(\cdot)$ be a sigmoidal function. Let I_m denote the m -dimensional unit hypercube $[0, 1]^m$. The space of continuous functions on I_m is denoted by $C(I_m)$.

Then, given any $\varepsilon > 0$ and any function $f \in C(I_m)$, there exist an integer N , real constants $v_i, b_i \in \mathbb{R}$ and real vectors $w_i \in \mathbb{R}^m$, where $i = 1, \dots, N$ such that we may define:

$$F(x) = \sum_{i=1}^N v_i \varphi \left(\sum_{k=1}^m w_{ik} x_k + b_i \right)$$

as an approximate realization of the function f , where f is independent of φ ; that is,

$$|F(x) - f(x)| < \varepsilon$$

for all $x \in I_m$. In other words, functions of the form $F(x)$ are dense in $C(I_m)$.



Neural Networks and Statistics



KTH Matematik

Let us write

$$C(\mathbf{w}, \mathbf{b}) = \frac{1}{2n} \sum_{r=1}^n \sum_{i=1}^K (y_i(\mathbf{x}^{(r)}) - t_i^{(r)})^2 = \frac{1}{2n} \sum_{r=1}^n \sum_{i=1}^K (t_i^{(r)} - y_i(\mathbf{x}^{(r)}))^2.$$

Set $(\mathbf{t}^{(r)}) = (t_1^{(r)}, \dots, t_K^{(r)})$. Let us assume that the pairs $(\mathbf{t}^{(r)}, \mathbf{x}^{(r)})$, $r = 1, \dots, n$ are independent samples with the probability distribution $P(\mathbf{T}, \mathbf{X})$. Then the Law of Large Numbers gives

$$C(\mathbf{w}, \mathbf{b}) \approx \frac{1}{2} E [\|\mathbf{T} - y(\mathbf{X})\|^2]$$

for large n , where

$$\|\mathbf{t} - y\|^2 = \sum_{i=1}^K (t_i - y_i)^2$$



$$C(\mathbf{w}, \mathbf{b}) \approx \frac{1}{2} E [\|\mathbf{T} - y(\mathbf{X})\|^2]$$

and probability theory tells that $E [\|\mathbf{T} - y(\mathbf{X})\|^2]$ is minimized by

$$y^*(\mathbf{X}) = E [\mathbf{T} | \mathbf{X}]$$

Hence we can say that a neural network after learning (with mean squared error cost criterion) approximates the conditional expectation of the target variable given the data.



For example, in a two-category classification problem, if the target t is coded as 1, if the object is from class 1 and -1 if it is from class 2. Then,

$$E[\mathbf{T} \mid \mathbf{X} = \mathbf{x}] = +1 \cdot P(\omega_1 \mid \mathbf{x}) + (-1) \cdot P(\omega_2 \mid \mathbf{x})$$

and the neural network estimates the following (well-known) discriminant function:

$$P(\omega_1 \mid \mathbf{x}) - P(\omega_2 \mid \mathbf{x}).$$



- M. D. Richard and R. Lippmann: Neural network classifiers estimate Bayesian a posteriori probabilities. *Neural Computation*, vol. 3, pp. 461– 483, 1991.



Bias -Variance Trade-Off

Suppose that we have a training set consisting of a set of points $\mathbf{x}_1, \dots, \mathbf{x}_n$ and real values t_i associated with each point \mathbf{x}_i . We assume that there is a functional, but noisy relation

$$t_i = f(\mathbf{x}_i) + \epsilon,$$

where the noise, ϵ , has zero mean and variance σ^2 .

We want to find a function $\hat{f}(\mathbf{x})$, that approximates the true function $t = f(\mathbf{x})$ as well as possible, by means of some learning algorithm. We make "as well as possible" precise by measuring the MSE between y and $\hat{f}(\mathbf{x})$ which we want to be minimal both for $\mathbf{x}_1, \dots, \mathbf{x}_n$ AND for points outside of our sample. Of course, we cannot hope to do so perfectly, since the y_i contain noise ϵ . This means we must be prepared to accept an irreducible error in any function we come up with.



Bias -Variance Decomposition

Finding an \hat{f} that generalizes to points outside of the training set can be done with any of the countless algorithms used for supervised learning. It turns out that whichever function \hat{f} we select, we can decompose its expected error on an unseen sample \mathbf{x} as follows:

$$\mathbb{E}[(t - \hat{f}(\mathbf{x}))^2] = \text{Bias}[\hat{f}(\mathbf{x})]^2 + \text{Var}[\hat{f}(\mathbf{x})] + \sigma^2 \quad (5)$$

(6)

Where:

$$\text{Bias}[\hat{f}(\mathbf{x})] = \mathbb{E}[\hat{f}(\mathbf{x})] - f(\mathbf{x}) \quad (7)$$

and

$$\text{Var}[\hat{f}(\mathbf{x})] = \mathbb{E}[(\hat{f}(\mathbf{x}) - \mathbb{E}[\hat{f}(\mathbf{x})])^2] \quad (8)$$



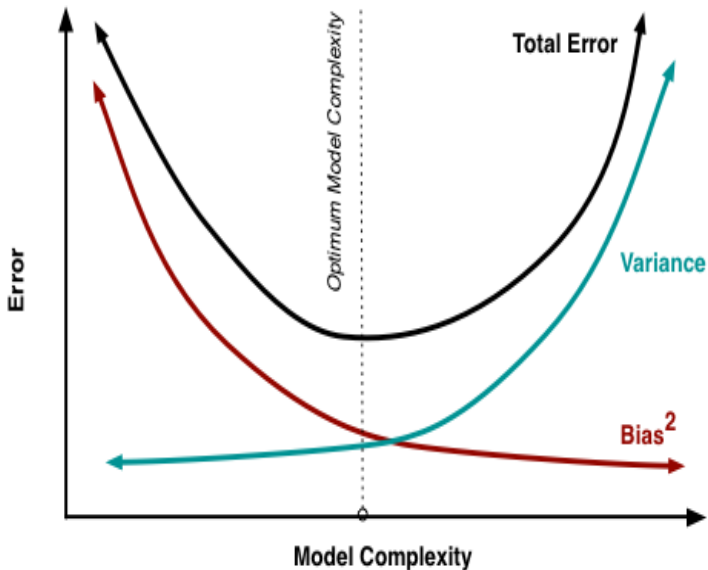
Bias -Variance Trade-Off

The expectation ranges over different choices of the training set $\mathbf{x}_1, \dots, \mathbf{x}_n, t_1, \dots, y_n$ all sampled from the same distribution. The three terms represent:

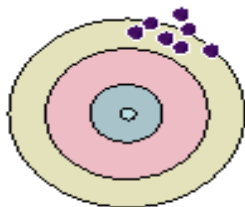
- the square of the bias of the learning method, which can be thought of the error caused by the simplifying assumptions built into the method. E.g., when approximating a non-linear function $f(\mathbf{x})$ using a learning method for linear models, there will be error in the estimates $\hat{f}(\mathbf{x})$ due to this assumption;
- the variance of the learning method, or, intuitively, how much the learning method $\hat{f}(\mathbf{x})$ will move around its mean;
- the irreducible error σ^2 . Since all three terms are non-negative, this forms a lower bound on the expected error on unseen samples.



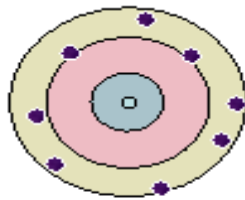
Bias -Variance Trade-Off



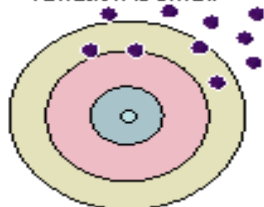
Bias -Variance Trade-Off



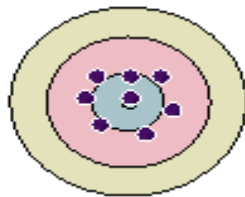
Bias is large
variation is small



Bias is small
variation is large



Bias is large
variation is large



Bias is small
variation is small

Accuracy versus Quality of an Estimator Using Bias and Variation as Measurable Quantities Respectively

Bias -Variance Trade-Off

- *If the model is too simple, the solution it yields is biased and does not fit the data.*
- *If the model is too complex, it is too sensitive to small variations in data*



Ideally, the optimal model that minimizes the overall MSE given by $E[\mathbf{T}|\mathbf{X}]$, which leaves the minimum MSE to be the intrinsic error σ^2 . Because of the randomness of the limited training set, the estimate ANN is also a random variable which will hardly be the $E[\mathbf{T}|\mathbf{X}]$ for a given training set. The bias and best possible function variance terms hence provide useful information on how the estimation differs from the desired function.



Bias -Variance Trade-Off

The model bias measures the extent to which the average of the estimation function over all possible data sets with the same size differs from the desired function. The model variance, on the other hand, measures the sensitivity of the estimation function to the training data set. Although it is desirable to have both low bias and low variance, we can not reduce both at the same time for a given data set because these goals are conflicting.

